

Актуальные угрозы ML-алгоритмов с точки зрения ИБ

Александра Мурзина, Positive Technologies



ML применяется в ИБ

- Для детектирования атак
 - Продукты
 - Экспертиза
- Как инструмент хакера
- Сам может быть уязвим

Есть ли проблема?

- Шифровальщики
- Сливы данных
- Ну дипфейки, но выглядит как шутка какая-то

Есть ли что-то, что приносит ущерб?

ML в нашей жизни

- Биометрия
- Рекомендательные системы
- Переводчики
- Системы анализа спама и фрода
- Голосовые ассистенты
- IoT, умные устройства

ML в нашей жизни

Можно ли представить жизнь без этих сервисов?

А как мы знаем, вопросы ИБ актуальны как никогда.

1. Какие атаки изучают?

2. Чего боятся вендоры?

3. Что происходит на самом деле?

Какие атаки изучают больше всего

Бывает, что все атаки на МЛ, которые есть, обобщают до Adversarial attack.

Какие атаки изучают больше всего

Бывает, что все атаки на МЛ, которые есть, обобщают до Adversarial attack.

Всегда найдутся данные, на которых алгоритм ошибется.
Задача злоумышленника – подобрать такие данные.
Как?

Какие атаки изучают больше всего

Бывает, что все атаки на МЛ, которые есть, обобщают до Adversarial attack.

Всегда найдутся данные, на которых алгоритм ошибется.
Задача злоумышленника – подобрать такие данные.
Как?

С помощью adversarial техник найти тот самый blind spot модели – такое место, где модель не уверена в принимаемом решении.

Adversarial attacks in the wild

Затруднение копирования контента

- Авито публиковали контент с замененным номерным знаком
- Авто.ру копировали, заменяя номерной знак на свой



<https://tinyurl.com/36348wd4>

Adversarial attacks in the wild

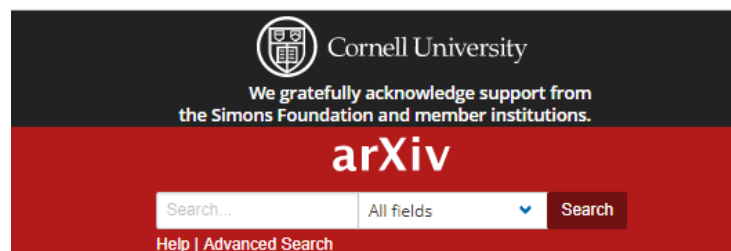
Затруднение копирования контента

- Авито публиковали контент с замененным номерным знаком
- Авто.ру копировали, заменяя номерной знак на свой
- Авито готовили изображение, используя Adversarial attack
- ML-модель Авто.ру больше не могла обнаруживать номерной знак и менять его на свой



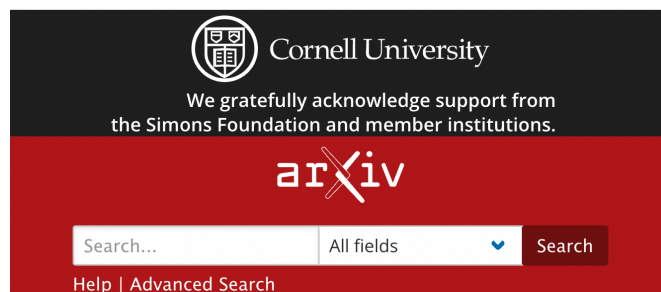
<https://tinyurl.com/36348wd4>

Какие атаки изучают больше всего



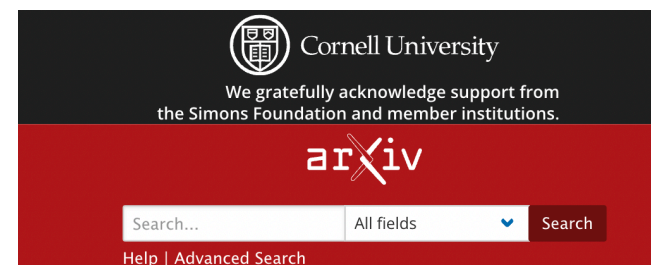
1,504

2019



5,322

апрель 2022



6,161

ноябрь 2022

Риски безопасности ИИ от Kang Li

AI App Security Risk

Model Security

- **Adversarial ML**
- Model Backdoor
- Model Theft

Implementation Security

- Sensor Security
- Flaws in Framework
- Logical Flaws

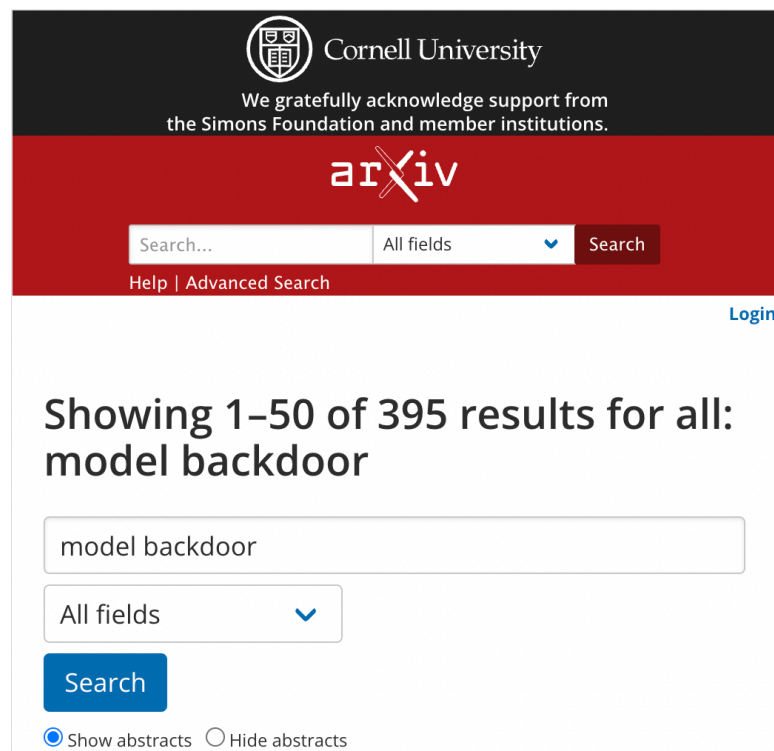
Data Integrity Security

- Data Poisoning
- Scaling Attack
- Risk over Network

<https://tinyurl.com/4fh7j3ky>

Какие атаки изучают больше всего

Model backdoors



The screenshot shows the arXiv search interface. At the top, there is a Cornell University logo and a message: "We gratefully acknowledge support from the Simons Foundation and member institutions." Below this is the arXiv logo. A search bar contains the text "Search..." and a dropdown menu is set to "All fields". A "Search" button is to the right of the search bar. Below the search bar, there are links for "Help" and "Advanced Search". On the right side of the header, there is a "Login" link. The main content area shows the search results: "Showing 1–50 of 395 results for all: model backdoor". Below this, there is a search bar containing the text "model backdoor", a dropdown menu set to "All fields", and a "Search" button. At the bottom, there are radio buttons for "Show abstracts" (which is selected) and "Hide abstracts".

- так какая атака опасней?
- та, которая причиняет наибольший ущерб
- так какая?
- а какие у вас риски?

Чего боятся вендоры? Обзор Gartner

Случались ли у вас нарушения конфиденциальности ИИ или инциденты безопасности?



41% опрошенных организаций столкнулись с нарушением конфиденциальности ИИ или инцидентом безопасности.

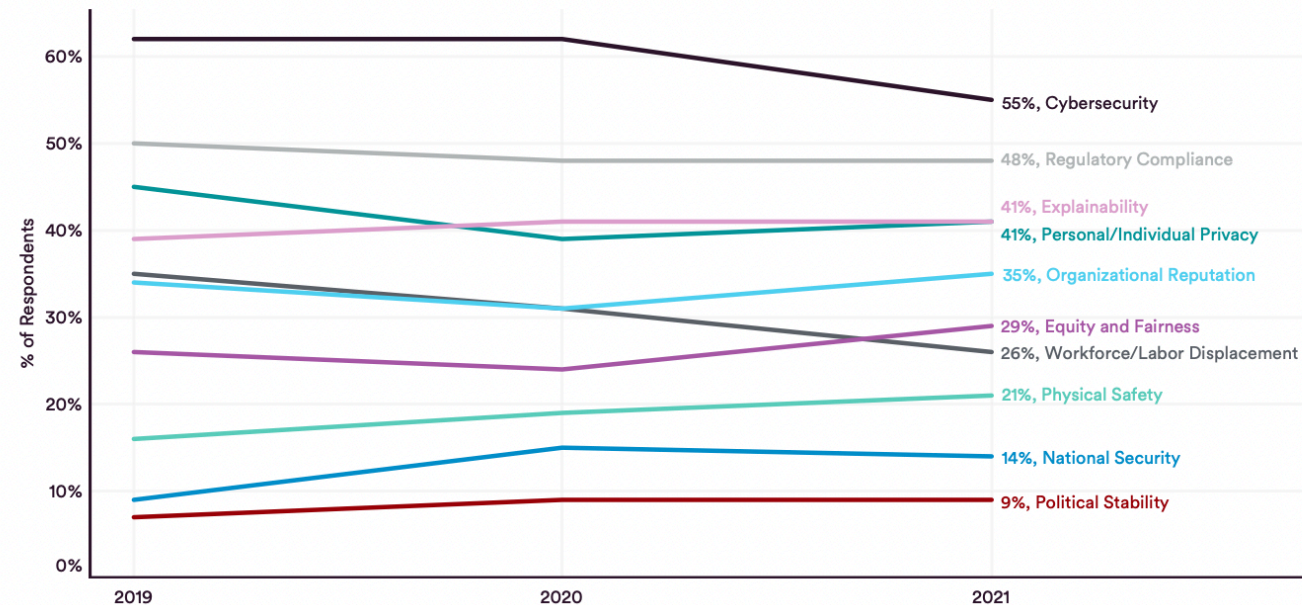
Из этих инцидентов 60% были компрометацией данных внутренней стороной, а 27% — злонамеренными атаками на ИИ-инфраструктуру организации.

<https://tinyurl.com/mr3fhkvx>

Есть ли проблема?

RISKS from ADOPTING AI that ORGANIZATIONS CONSIDER RELEVANT, 2019–21

Source: McKinsey & Company, 2021 | Chart: 2022 AI Index Report



Наиболее актуальным риском в 2021 году является кибербезопасность (55% респондентов).

Далее следуют нормативно-правовое соответствие (48%), интерпретируемость (41%) и конфиденциальность (41%).

Consideration and Mitigation of Risks From Adopting AI

<https://tinyurl.com/3kax9hbw>

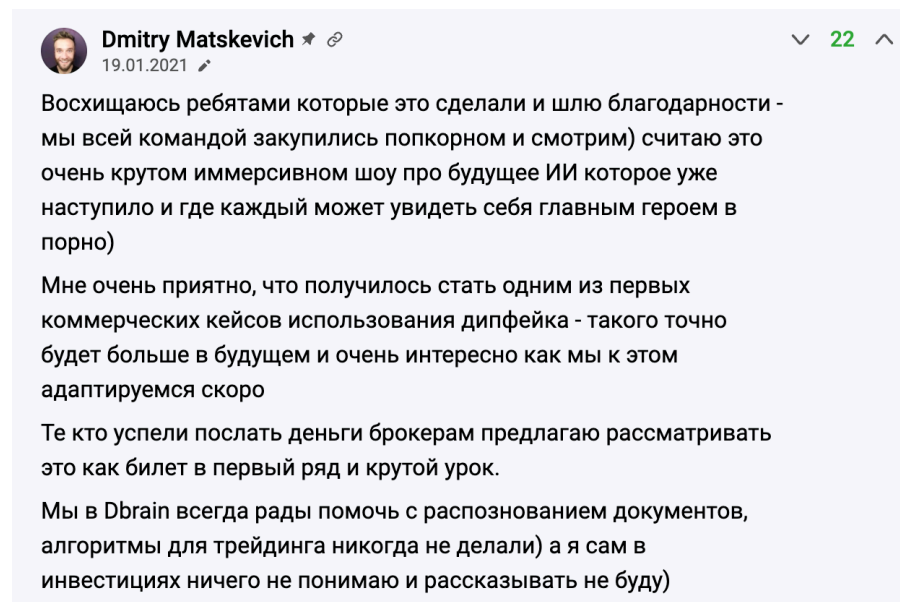
Что происходит на самом деле?

спойлер:

нет официальной статистики по инцидентам

Сделали копию лица для DeepFake-видео

Неизвестные скопировали лицо основателя Dbrain
Дмитрия Мацкевича для deepfake-видео с рекламой супердоходов.



<https://tinyurl.com/2x7755vw>

Оформить кредит в банке по биометрии?

29 марта 2021

«Компания расширит функции запущенного в 2018 году приложения «Биометрия» услугой дистанционной сдачи биометрии. Это предложение рассматривается ФСБ и Федеральной службой по техническому и экспортному контролю (ФСТЭК)»

<https://tinyurl.com/yepaz9r9>

Оформить кредит в банке по биометрии?

29 марта 2021

«Компания расширит функции запущенного в 2018 году приложения «Биометрия» услугой дистанционной сдачи биометрии. Это предложение рассматривается ФСБ и Федеральной службой по техническому и экспортному контролю (ФСТЭК)»

<https://tinyurl.com/yepaz9r9>

10 апреля 2021

В Москве мошенники стали использовать голоса клиентов банков для оформления кредитов.

<https://tinyurl.com/ypyhje2y>

Можно ли манипулировать ML-системой?

Полиция в США включает популярную музыку на забастовках, чтобы исключить попадание видео с собой в соц сети.

“

You can record all you want, I just know it can't be posted to YouTube.

<https://tinyurl.com/5b9suxjy>

А нас это как касается?

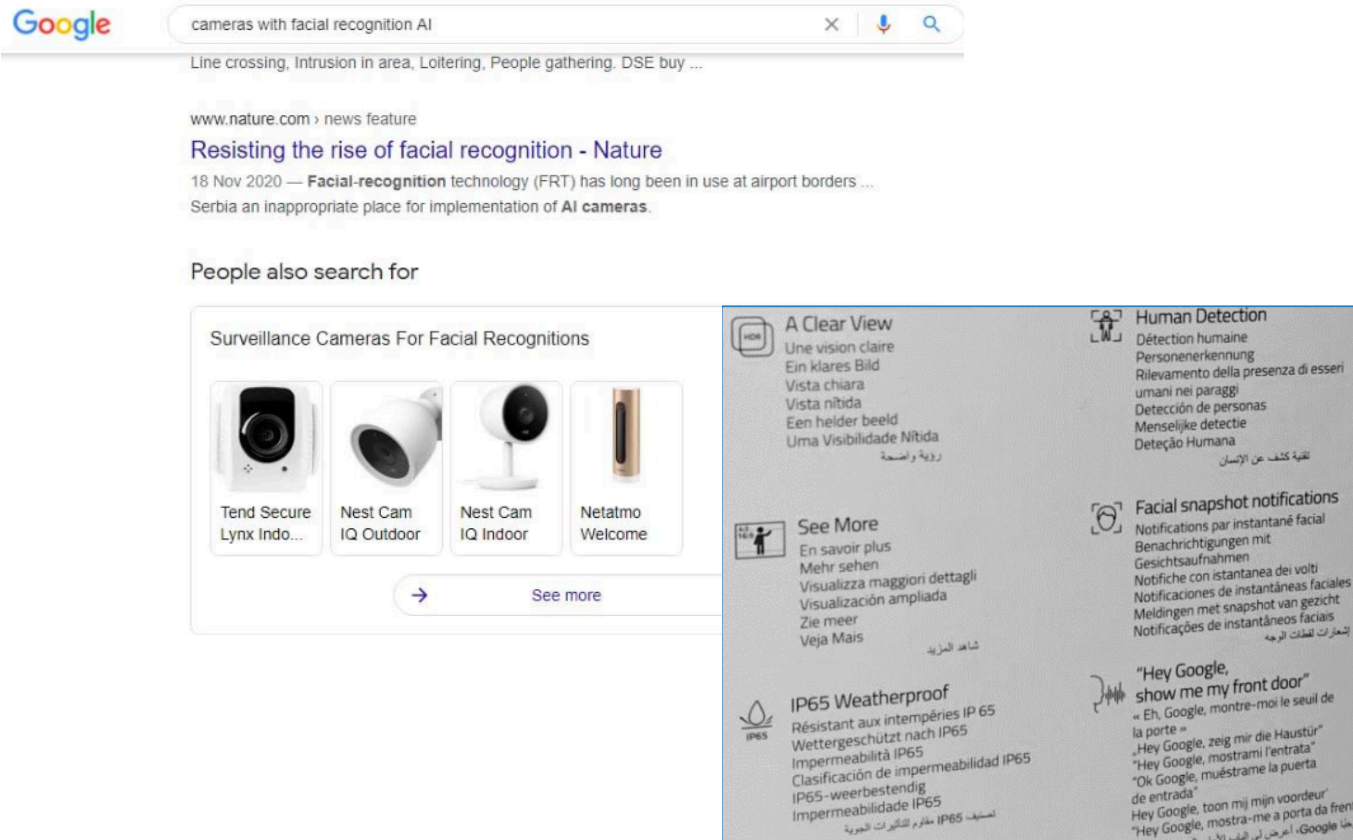


IoT & биометрия

Исследования Positive Technologies:

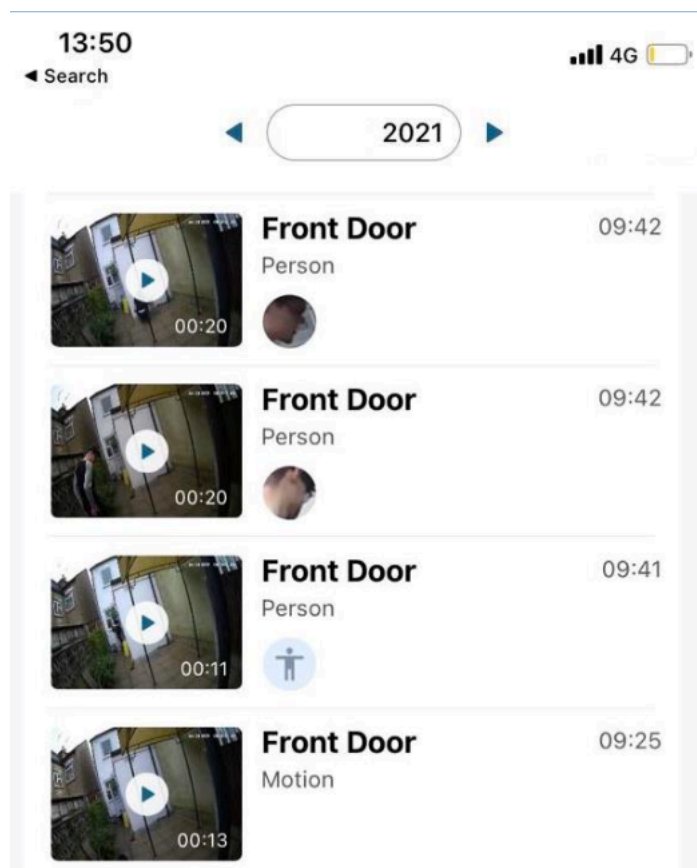
- Умный дверной звонок
- Устройство биометрической идентификации

Умный дверной звонок



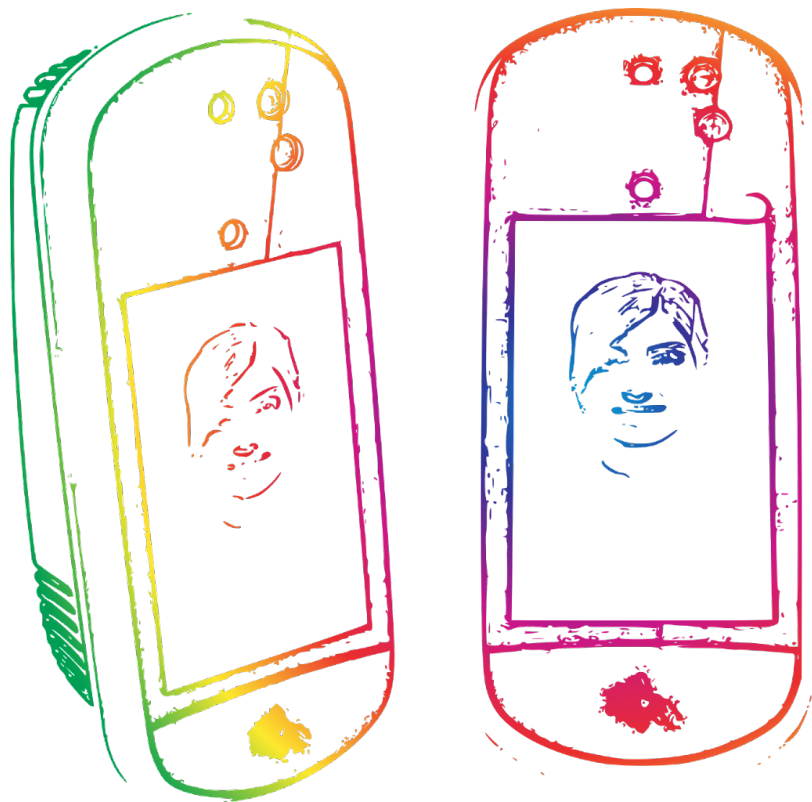
- Вендор заявляет о безопасности.
- Может открывать замок автоматически.
- Может быть подключен к общему хабу умного дома.

Умный дверной звонок



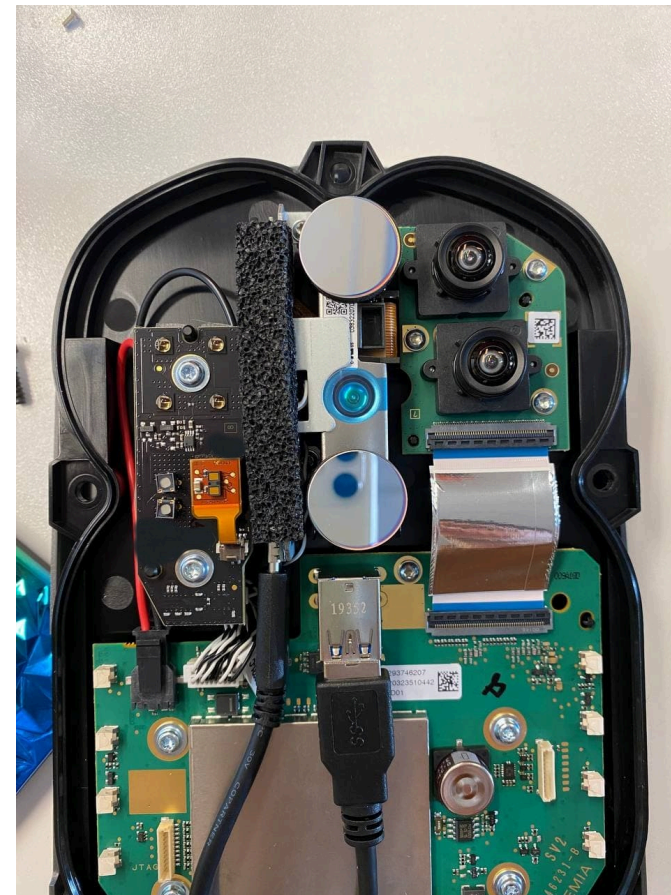
- Вендор заявляет о безопасности.
- Может открывать замок автоматически.
- Может быть подключен к общему хабу умного дома.
- Слабое шифрование на этапе авторизации.
- Возможность получить доступ к хабу умного дома.
- Можно следить за тем, кто приходит/уходит.

Умное устройство биометрической идентификации



- Камера глубины
- 2 камеры видимого диапазона

<https://youtu.be/2fvdyKl1NiY>



Умное устройство биометрической идентификации

Как работает?

1. Детектируем лицо в кадре.
2. Проверяем Liveness камерой глубины.
3. Захватываем лицо из камеры видимого диапазона.
4. Предобработка.
5. DNN.
6. Сравнение с БД.

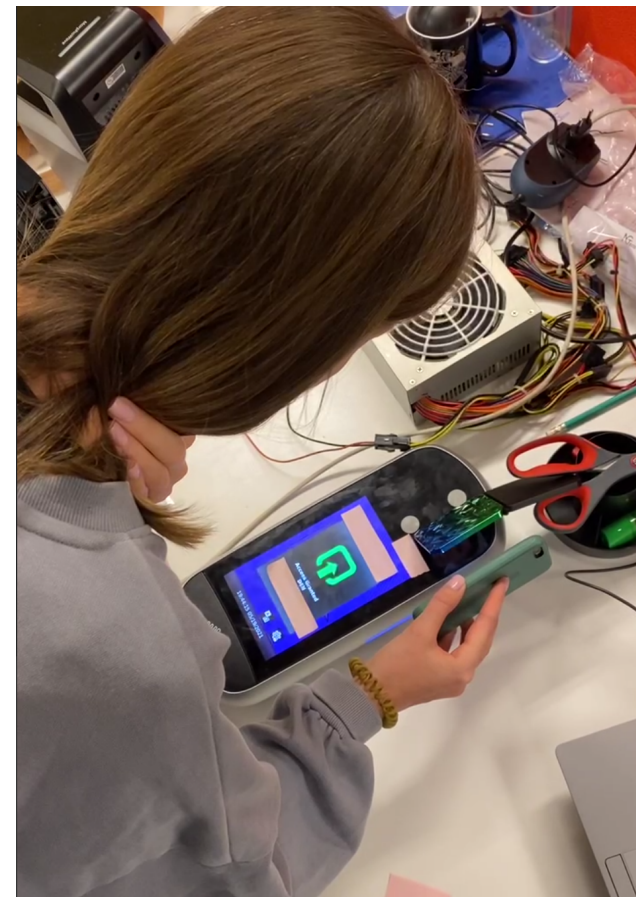
<https://youtu.be/2fvdyKl1NiY>

Умное устройство биометрической идентификации

Как работает?

1. Детектируем лицо в кадре.
- 2. Проверяем Liveness камерой глубины.**
- 3. Захватываем лицо из камеры видимого диапазона.**
4. Предобработка.
5. DNN.
6. Сравнение с БД.

<https://youtu.be/2fvdyKl1NiY>



KYC

KYC (Know Your Customer) — обязательная для финансовых институтов процедура идентификации контрагентов. Включает в себя процедуры отбора и идентификации, а также отслеживание транзакций и их анализ.

- Онлайн — наше всё.
- Но надо как-то подтверждать свою личность и далее транзакции.

<https://tinyurl.com/2s4vhkvj>



Что делать?

Зарегулировать!

Совфед запретит искусственному
интеллекту дискриминировать
россиян

**The EU wants to put
companies on the hook
for harmful AI**

<https://tinyurl.com/ykc56tkr>

<https://tinyurl.com/3h7526hz>

Может, не надо?

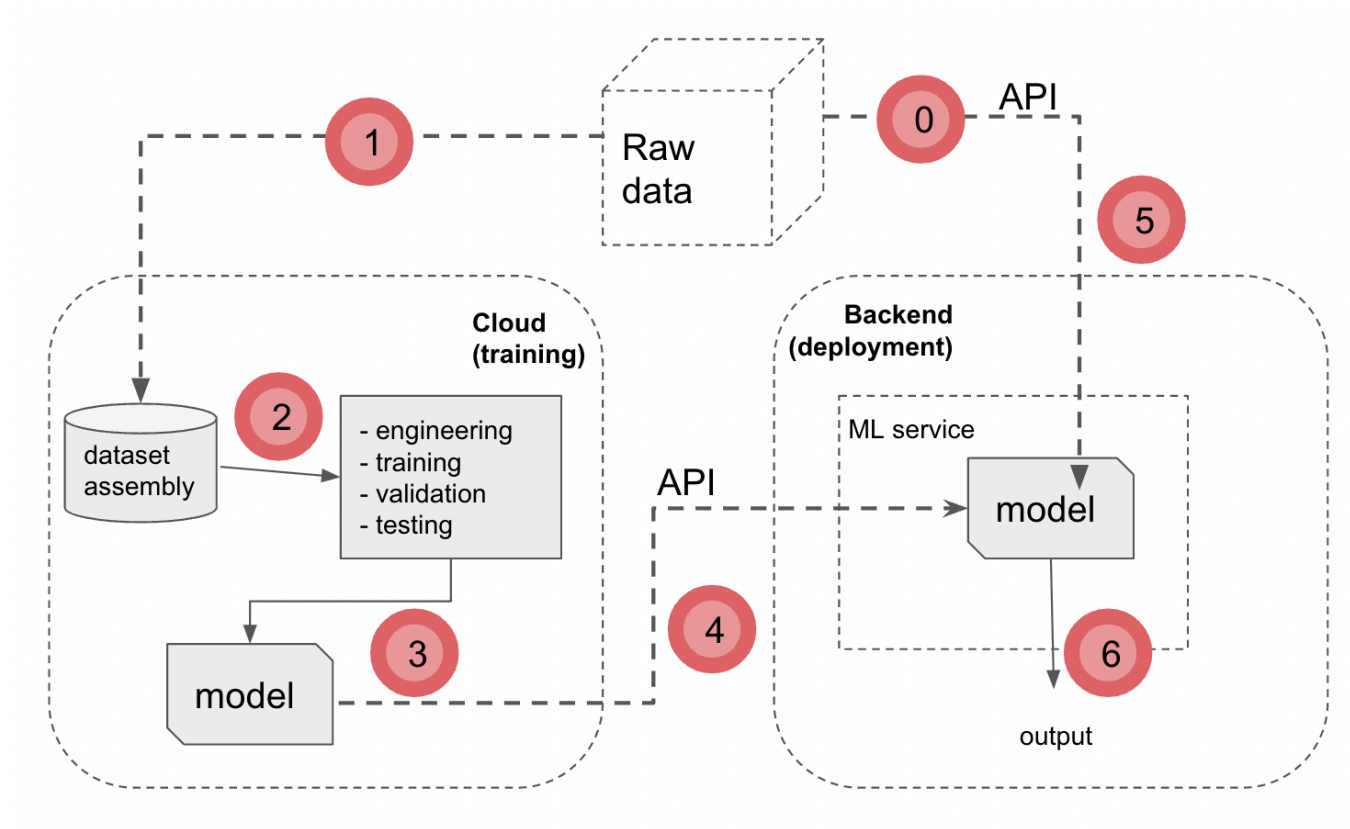
George Hotz, основатель стартапа Comma.AI:

Today's struggle isn't a technological one, it's psychosocial. The technology of sovereignty is only part of the battle, particularly if that technology is further down the tech tree than non invasive wireheading.

Comma.AI — опенсорс-ассистент водителя с поддержкой 200+ машин — опенсорс-беспилотник.

<https://tinyurl.com/ccczpfd8>

Модель угроз



- 0 — api security
- 1 — poisoning
- 2 — trojanning
- 3 — backdooring
- 4 — model extraction
- 5 — model inversion
- 6 — inference attacks

И что с этим делать?

- Багбаунти
- Киберполигоны
- Внешний аудит

<https://tinyurl.com/mut5d6z7>

<https://tinyurl.com/42vcnpzp>

Lessons learned

- То, что изучают, - не всегда то, чего боятся вендоры и не всегда то что в происходит в реальной жизни.
- Важно понимание угроз, которые кажутся наиболее опасными.
- Стоит начать с классической безопасности.
- Некоторых рисков ИБ в МЛ можно избежать, если учесть в архитектуре решения.
- Это бесконечная игра.

Актуальные угрозы ML-алгоритмов с точки зрения ИБ

Александра Мурзина, Positive Technologies

@murzina_a

за доклад можно
голосовать тут

